

# DOCUMENTS CLUSTERING BASED ON MAX-CORRENTROPY NONNEGATIVE MATRIX FACTORIZATION

Le Li<sup>1</sup>, Jianjun Yang<sup>2</sup>, Yang Xu<sup>3</sup>, Zhen Qin<sup>3</sup>, Honggang Zhang<sup>3</sup>

<sup>1</sup>David R. Cheriton School of Computer Science, University of Waterloo, ON N2L3G1, Canada

<sup>2</sup>Department of Computer Science, University of North Georgia, Oakwood, GA 30566, USA

<sup>3</sup>Pattern Recognition and Intelligent System Laboratory, Beijing University of Posts and Telecommunications, Beijing, China  
E-MAIL: l248li@uwaterloo.ca, jianjun.yang@ung.edu, {xj992adolphy, qinzhenbupt}@gmail.com, zhhg@bupt.edu.cn

## Abstract:

Nonnegative matrix factorization (NMF) has been successfully applied to many areas for classification and clustering. Commonly-used NMF algorithms mainly target on minimizing the  $l_2$  distance or Kullback-Leibler (KL) divergence, which may not be suitable for nonlinear case. In this paper, we propose a new decomposition method by maximizing the correntropy between the original and the product of two low-rank matrices for document clustering. This method also allows us to learn the new basis vectors of the semantic feature space from the data. To our knowledge, we haven't seen any work has been done by maximizing correntropy in NMF to cluster high dimensional document data. Our experiment results show the supremacy of our proposed method over other variants of NMF algorithm on Reuters21578 and TDT2 datasets.

## Keywords:

Document clustering; Nonnegative matrix factorization

## 1. Introduction

A corpus is a collection of documents where each document is associated with a ground-truth topic that summarizes the content of the document. Document clustering is the process that finds the correct label for the input document, such that this label should match with the ground-truth topic as much as possible. Such clustering makes it possible that automatically organizes millions of documents, websites, news, etc. into the multiple partitions, where documents within the same partitions share same topic. As a consequence, we can leverage this technique to different tasks, like document organization and browsing, corpus summarization, and document classification [1].

Different types of algorithms have been used to cluster/classify

the data (e.g. SVM [33] and pLSA [32]). These algorithms have a variety of applications in different areas [30, 22, 18, 26, 31, 20, 19, 17, 16, 25]. Among those algorithms, we are especially interested in the nonnegative matrix factorization (NMF) method. NMF algorithm maps the original features into latent semantics space where each basis vector in the latent space represents a topic. More precisely, assuming each document is represented as a feature vector with  $D$  dimension, and we have  $N$  documents in the corpora, then we can form a  $D * N$  matrix (denoted as  $X$ ) to represent the whole corpora. NMF algorithm can decompose the  $X$  into two low-rank nonnegative matrices,  $H$  and  $W$ , such that the  $X \approx HW$ . One of the main benefits is its nature of dimension reduction without losing too much useful information. This decomposition has been shown its supremacy in many areas (e.g. bioinformatics [23]).

Much work has been done on applying NMF algorithms to document clustering [27, 15]. However, most of them try to minimize the  $l_2$  distance or KL divergence. Inspired by the recent work in [23] that combines correntropy with NMF in cancer clustering, we propose a similar max-correntropy nonnegative matrix factorization algorithm (MCC) into the document clustering area. The work in [23] is in line with ours in the way that both show the benefits of using this max-correntropy method for clustering. However, we are working on different areas. Meanwhile, the work in [23] only examines the clustering performance on a limit number of topics (less than 10) and lower dimension data, while we systematically investigate its performance on more sophisticated clustering tasks with more documents, topics and higher dimension of data.

To achieve that, we implement the MCC algorithm and test its accuracy on the Reuters21578 and TDT2 corpora. We compare the MCC algorithm to classic loss functions ( $l_2$  distance and KL divergence), as well as other variants of NMF

algorithms. The results show that the proposed algorithm suppresses the rest methods on document clustering in both datasets. Moreover, we fully investigate how the MCC algorithm behaves when we keep increasing the number of topics (i.e. increasing the clustering difficulty). We find that although all algorithms' accuracies drop as we increase the number of topics, MCC algorithm is the most robust algorithm against the increment of the number of topics.

The main contribution made in this paper is that, to our knowledge, this is the first work that factorizes the matrix from the perspectives of maximizing correntropy in document clustering. Besides that, we fully investigate its behaviors when faced with multiple topics' documents and high-dimension data. The results show the benefits of using the proposed MCC algorithm in document clustering.

## 2 Related work

$K$ -means is a classic clustering algorithm. This algorithm finds the closest cluster for each document by finding the smallest distance between the document and the existing clusters' centroids. The clusters' centroids are also updated at each iteration due to new cluster members.  $K$ -means is based on the assumption that documents belonging to same topic should also be close to each other in the feature space. In a similar vein, Naive Bayes and Gaussian mixture model [14, 3, 13] are used based on different document distribution assumptions. One problem with these methods is that if the corpora properties don't following such assumptions, the performance of these algorithms may be at risk.

Latent Semantics Indexing (LSI) [6] is the technique that converts the corpora from the original feature space into a latent semantics space. Each basis axis in the latent semantics space essentially represents one type of semantic information of the corpora. By doing so, each document is essentially a combination of multiple semantics information. Then we can apply the classic clustering algorithms on these new representations of documents in latent semantics space. One issue with this method is that the coefficients of the combination could be positive or negative. A negative coefficient is not such a natural way to interpret the document. Meanwhile, the bases that spanning the latent semantics space in some LSI algorithms, like Singular Vector Decomposition (SVD) [7], are orthogonal, which means that every semantics bases are different from each other. However, in reality, this is not always the case.

Similar to LSI algorithms, NMF also maps the corpora into latent feature space. The differences are that: firstly, the bases in latent feature space don't need to be orthogonal. Also, each

basis now corresponds to one topic of the corpora, which makes it very easy to determine the topic of document by simply choosing the largest component in the latent space. Meanwhile, every element in the two decomposed low-rank matrices are nonnegative. This additive combination makes it more natural to understand each document in an intuitive manner.

The benefits of using NMF in document clustering have been heavily investigated in many existed papers [27, 15, 12, 24]. However, many of them are mainly targeting on minimizing the  $l_2$  norm or KL divergence in the process of matrix decomposition. Correntropy-based decomposition methods have proved effective in many areas like cancer clustering [23], face recognition [9], etc. Some other solutions can be found in [21, 4]. However, we never find such technique in the document clustering research or be used for very high-dimension data with considerable number of clusters, which is another starting point of our work.

## 3 Algorithm

Assuming we have a matrix  $X \in \mathbb{R}^{D \times N}$ . NMF allows us to factorize  $X$  into two nonnegative matrices  $H \in \mathbb{R}^{D \times K}$  and  $W \in \mathbb{R}^{K \times N}$ , where the product  $H * W$  approximates the original matrix  $X$ . Each column in  $X$  is the feature vector of one document with  $D$  elements. Thus,  $X$  essentially represents the whole corpus with  $N$  documents. Conventionally, we name  $H$  as basis matrix that each column forms the basis vector of the semantic feature space, and  $W$  as coefficient matrix. Hence, a document is further represented as the additive combination of weighted basis vectors in semantic space.  $l_2$  norm and Kullback-Leibler (KL) divergence are two commonly-used measures of the similarity between original matrix  $X$  and the product of  $H$  and  $W$ . Based on different similarity measures, we are able to solve the factorization problem by minimizing the corresponding errors between  $X$  and  $H * W$ .

In this paper, we propose a new method to quantify the NMF by maximizing correntropy criteria in document clustering. Correntropy measures the generalized similarity between two random variables. More precisely, it models the expected differences between two random variables after mapping through kernel function. Without knowing the joint distribution of  $X$  and  $Y$ , we can simply estimate the expectation by taking average (shown in Equation 1):

$$\hat{V}_\sigma(x, y) = \frac{1}{D} \sum_{i=1}^D k_\sigma(x_i - y_i) \quad (1)$$

where  $k_\sigma(\cdot)$  is the kernel function and  $x_i$  and  $y_i$  are the element in  $X$  and  $Y$ , respectively.

Thus, instead of using  $l_2$  distance or KL divergence, we try to find the basis matrix  $H$  and coefficient matrix  $W$ , whose product  $Y$  approximates  $X$ , by maximizing their correntropy on a feature-by-feature basis to allow for weighting each feature differently. For each feature, the kernel function can be calculated as:

$$k_\sigma \left( \sqrt{\sum_{n=1}^N (x_{dn} - \sum_{k=1}^K h_{dk} w_{kn})^2} \right) \quad (2)$$

Hence, the correntropy maximization problem is expressed in Equation 3.

$$\max_{h_{dk} > 0, w_{kn} > 0} \sum_{d=1}^D k_\sigma \left( \sqrt{\sum_{n=1}^N (x_{dn} - \sum_{k=1}^K h_{dk} w_{kn})^2} \right) \quad (3)$$

To simplify the calculations without losing generality, we choose the Gaussian kernel function as  $k_\sigma(\cdot)$ :

$$k_\sigma(x - y) = \exp(-\gamma \|x - y\|^2) \quad (4)$$

After substituting Equation 4 back into Equation 3, the basis and coefficient matrices can be derived by solving:

$$\max_{h_{dk} > 0, w_{kn} > 0} \sum_{d=1}^D \exp \left( -\gamma \sum_{n=1}^N (x_{dn} - \sum_{k=1}^K h_{dk} w_{kn})^2 \right) \quad (5)$$

We introduce the convex conjugate function  $\varphi(\cdot)$  and auxiliary variables  $\rho = [\rho_1, \dots, \rho_D]^\top$ . Based on the theory of convex conjugate functions, the above optimization problem is equivalent to:

$$\begin{aligned} & \max_{H, W, \rho} \hat{F}(H, W, \rho) \\ & \text{s.t. } H \geq 0, W \geq 0 \end{aligned} \quad (6)$$

$$\hat{F}(H, W, \rho) = \sum_{d=1}^D \left( \rho_d \sum_{n=1}^N (x_{dn} - \sum_{k=1}^K h_{dk} w_{kn})^2 - \varphi(\rho_d) \right)$$

The optimization problem can be solved by Expectation-Maximization-like method. Starting from the initial value of  $H$  and  $W$ , we compute  $\rho$  in expectation step (E-step). Conditional on the  $\rho$  value, we update the  $H$  and  $W$  values in maximization step (M-step). The process is called one iteration. This iterative process stops until it converges. The proposed MCC has a good convergence performance. We direct the readers to refer similar convergence proof in [23]. We often assign  $H$  and  $W$  with random numbers to start the algorithm if we have no prior information about the distribution of data.

**E-step:** Starting from the estimated  $H$  and  $W$  from last M-step (or random values in the 1st iteration),  $\rho$  of the  $t$ -th iteration is computed as:

$$\rho_d^t = -g \left( \sqrt{\sum_{n=1}^N \left( x_{dn} - \sum_{k=1}^K h_{dk}^t w_{kn}^t \right)^2}, \sigma^t \right) \quad (7)$$

$$\text{where } g(z, \sigma) = \sup_{\varrho \in \mathbb{R}^-} \left( \varrho \frac{\|z\|^2}{\sigma^2} - \varphi(\varrho) \right)$$

$$\text{and } \sigma^t = \sqrt{\frac{\theta}{2D} \sum_{d=1}^D \sum_{n=1}^N \left( x_{dn} - \sum_{k=1}^K h_{dk}^t w_{kn}^t \right)^2}$$

**M-step:** Conditional on the new  $\rho$  from last step, we compute the new basis and coefficient matrix, denoted as  $H^{t+1}$  and  $W^{t+1}$  respectively, by maximizing the object function:

$$\begin{aligned} & (H^{t+1}, W^{t+1}) \\ &= \arg \max_{H, W} \sum_{d=1}^D \left( \rho_d^t \sum_{n=1}^N \left( x_{dn} - \sum_{k=1}^K h_{dk} w_{kn} \right)^2 \right) \\ &= \arg \max_{H, W} \text{Tr}[(X - HW)^\top \text{diag}(\rho^t)(X - HW)] \\ &= \text{Tr}[X^\top \text{diag}(-\rho^t)X] - 2\text{Tr}[X^\top \text{diag}(-\rho^t)HW] \\ &\quad + \text{Tr}[W^\top H^\top \text{diag}(-\rho^t)HW] \\ &\text{s.t. } H \geq 0, W \geq 0 \end{aligned}$$

where  $\text{Tr}(\cdot)$  means the trace of the matrix.

We apply the Lagrange method to solve the optimization problem. Let the elements of matrices  $\Phi = [\phi_{dk}]$  and  $\Psi = [\psi_{kn}]$  be the corresponding Lagrange multipliers for the non-negative conditions of  $h_{dk} \geq 0$  and  $w_{kn} \geq 0$ . Then we can express the Lagrange optimization problem as:

$$\begin{aligned} L = & \text{Trac}[X^\top \text{diag}(-\rho^t)X] - 2\text{Trac}[X^\top \text{diag}(-\rho^t)HW] \\ & + \text{Trac}[W^\top H^\top \text{diag}(-\rho^t)HW] + \text{Trac}[\Phi H^\top] \\ & + \text{Trac}[\Psi H^\top] \end{aligned} \quad (8)$$

The partial derivatives of  $L$  with respect to  $H$  and  $W$  are:

$$\begin{aligned} \frac{\partial L}{\partial H} &= -2\text{diag}(-\rho^t)XW^\top + 2\text{diag}(-\rho^t)HWW^\top + \Phi \\ \frac{\partial L}{\partial W} &= -2H^\top \text{diag}(-\rho^t)X + 2H^\top \text{diag}(-\rho^t)HW + \Psi \end{aligned}$$

Set them to 0 based on Karush-Kuhn-Tucker optimal conditions, we have:

$$\begin{aligned} & -(\text{diag}(-\rho^t)XW^\top)_{dk}h_{dk} + (\text{diag}(-\rho^t)HWW^\top)_{dk}h_{dk} = 0 \\ & -(H^\top \text{diag}(-\rho^t)X)_{kn}w_{kn} + (H^\top \text{diag}(-\rho^t)HW)_{kn}h_{kn} = 0 \end{aligned}$$

Hence, the basis matrix  $H$  and coefficient matrix  $W$  can be updated as shown in Equation 9 and Equation 10.

$$h_{dk}^{t+1} \leftarrow h_{dk}^t \frac{(\text{diag}(-\rho^t)XW^{t\top})_{dk}}{(\text{diag}(-\rho^t)H^tW^tW_{dk}^{t\top})_{dk}} \quad (9)$$

$$w_{kn}^{t+1} \leftarrow w_{kn}^t \frac{(H^{t+1\top} \text{diag}(-\rho^t)X)_{kn}}{(H^{t+1\top} \text{diag}(-\rho^t)H^{t+1}W^t)_{kn}} \quad (10)$$

## 4 Experiment settings

We test the MCC algorithm on two datasets: Reuters21578<sup>1</sup> and TDT2<sup>2</sup>. These two datasets have been widely used in many places [27, 15] for document clustering. Reuters21578 test collection contains 21578 documents from 135 topics in total. We exclude those documents that belong to more than 1 topics in our experiment since we are trying to cluster each document to one single topic. Meanwhile, we also exclude those topics with less than 5 documents. As a consequence, we use 9545 documents for 51 topics in our experiment. TDT2 dataset contains around 11201 documents for 96 topics. We also apply similar pre-processing ways to it. The largest 30 topics with 9394 documents are used.

We use tf-idf method to extract the feature for each document. Stopwords and stemming are applied. The number of elements for each document are 16777 and 36771 for Reuters21578 and TDT2, respectively.

After matrix decomposition, the matrix  $W$  with dimension  $K \times N$  is essentially the new representation of the corpora in the way that each column is the feature vector of one document after dimension deduction. And the new dimension of the feature vector is  $K$  now. To evaluate the decomposition performance, we directly apply  $K$ -means clustering method to cluster  $W$  into  $K$  clusters.  $K$ -means will assign each document with a label. We compare the label from  $K$ -means to the original ground-truth label to calculate the clustering accuracy. The accuracy is defined in Equation 11.

$$\text{Accuracy} = \sum_{i=1}^N \delta(\text{kmeans\_label}_i, \text{topic}_i) / N \quad (11)$$

where  $\delta(\text{kmeans\_label}_i, \text{topic}_i)$  is the delta function, which returns 1 if  $\text{kmeans\_label}_i = \text{topic}_i$ ; otherwise 0. To find the correspondence between the topic from ground-truth data and the label by  $K$ -means, we use Kuhn-Munkres algorithm [10].

## 5 Results

One important parameter we need to control is the number of cluster  $K$ . Intuitively, the value of  $K$  controls the way to decompose the matrix. More importantly, it determines the number of topics that NMF algorithm can handle. That's to say, if  $K = 2$ , then the NMF is a two-topic clustering problem. If  $K > 2$ , then it's a multi-topic clustering problem. To fully investigate the efficiency of the MCC, we use the following candidate numbers of clusters: {2-10, 20, 30, 40, 51} for Reuters dataset and {2-10, 15, 20, 25, 30} for TDT2 corpus. We randomly select  $K$  topics, and use all documents from those  $K$  topics as the current run's testing documents. We repeat this process 20 times to reduce the potential effect of random errors on our experiment results since the performance of NMF algorithm is affected by the initial values of the iterative process. The average of 20 runs is used as the output of each algorithm.

Table 1 summarizes the results of two datasets with less than or equal to 10 topics. Figure 1 and Figure 2 demonstrate the accuracies on all chosen numbers of clusters. We firstly test the proposed MCC algorithm against two classic loss functions:  $l_2$  distance and KL divergence. It's clear that MCC algorithm outperforms the  $l_2$  distance and KL divergence in all cases of  $K$  in two datasets. This shows the supremacy of the MCC algorithm against the others. One possible reason is that  $l_2$  and KL distance are effective when dealing with linear separable data. However, if the data distribution is nonlinear manifold, it is considerably difficult for these two linear kernels to distinguish them.

Meanwhile, we observe that for all algorithms, the accuracy decreases as the number of clusters increases. Intuitively, more clusters inevitably increase the difficulties of finding the right label for each document. However, MCC is more robust to the increment of  $K$ , compared to other distance functions.

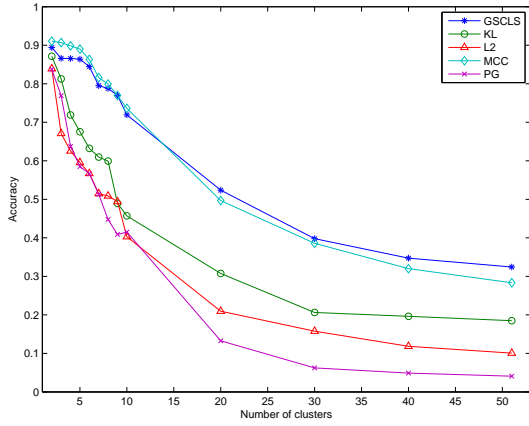
We also compare MCC against two variants of NMF algorithms: gradient descent-constrained least squares (GSCLS) [15], and Projected Gradient nonnegative matrix factorization (PG) [12]. Based on the results of two datasets, we can see that MCC suppresses the rest NMF algorithms when the number of clusters is smaller or equals to 10 on Reuters21578. When it comes to TDT2 dataset, MCC achieves the best performance in all cases, which shows the benefit of introducing the coreentropy into the factorization process. One potential reason is that MCC can self-learn different kernels for different features. This adaptive learning property somehow further improves the performance of MCC when facing with nonlinear datasets (e.g. document collection).

<sup>1</sup><http://kdd.ics.uci.edu/databases/reuters21578/reuters21578.html>

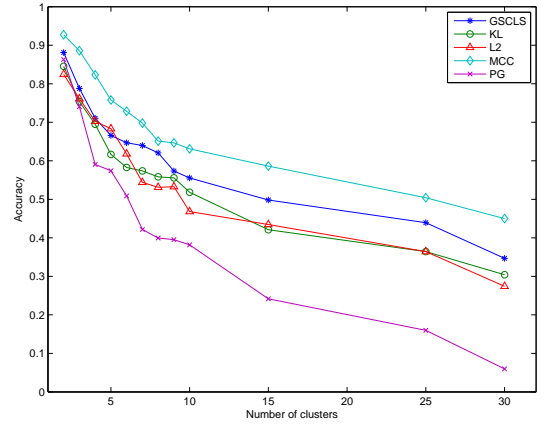
<sup>2</sup><http://www.itl.nist.gov/iad/mig/tests/tdt/1998/>

**Table 1. Clustering accuracies on Reuters21578 and TDT2 datasets**

| Number of clusters | Reuters21578 |       |       |        |       | TDT2  |       |       |        |       |
|--------------------|--------------|-------|-------|--------|-------|-------|-------|-------|--------|-------|
|                    | MCC          | L2    | K-L   | GS-CLS | PG    | MCC   | L2    | K-L   | GS-CLS | PG    |
| 2                  | 0.911        | 0.839 | 0.871 | 0.894  | 0.838 | 0.927 | 0.824 | 0.845 | 0.881  | 0.863 |
| 3                  | 0.907        | 0.671 | 0.813 | 0.866  | 0.769 | 0.886 | 0.761 | 0.754 | 0.788  | 0.740 |
| 4                  | 0.898        | 0.625 | 0.719 | 0.866  | 0.637 | 0.823 | 0.702 | 0.695 | 0.710  | 0.591 |
| 5                  | 0.890        | 0.596 | 0.675 | 0.864  | 0.585 | 0.758 | 0.683 | 0.616 | 0.666  | 0.574 |
| 6                  | 0.863        | 0.567 | 0.632 | 0.845  | 0.567 | 0.729 | 0.618 | 0.583 | 0.647  | 0.509 |
| 7                  | 0.816        | 0.515 | 0.610 | 0.795  | 0.513 | 0.698 | 0.544 | 0.574 | 0.640  | 0.422 |
| 8                  | 0.799        | 0.509 | 0.599 | 0.788  | 0.448 | 0.651 | 0.531 | 0.558 | 0.621  | 0.400 |
| 9                  | 0.770        | 0.494 | 0.490 | 0.770  | 0.409 | 0.647 | 0.533 | 0.556 | 0.573  | 0.395 |
| 10                 | 0.736        | 0.403 | 0.457 | 0.719  | 0.415 | 0.631 | 0.468 | 0.518 | 0.555  | 0.382 |



**Figure 1. Accuracies on Reuters21578.**



**Figure 2. Accuracies on TDT2.**

## 6 Conclusion

In this paper, we propose a new method to decompose the matrix into two low-rank matrices by maximizing the correntropy between them, such that we can easily and effectively use the decomposed matrix to cluster high-dimension data. We test the proposed MCC algorithm in the application of document clustering. We compare our proposed method to other loss functions and NMF algorithms. The results demonstrate the supremacy of our method on Reuters21578 and TDT2 corpora in terms of accuracy. In future, we will investigate the possibility of our proposed method in medical instrument[2], mechanical instrument[5] and other related areas [29, 28, 8]. [32, 11]

## Acknowledgements

This work was partially supported by National Natural Science Foundation of China under Grant No.61273217, 61175011 and 61171193, the 111 project under Grant No.B08004.

## References

- [1] C. C. Aggarwal and C. Zhai. A survey of text clustering algorithms. In *Mining Text Data*, pages 77–128. Springer, 2012.
- [2] A. Anderson, P. K. Douglas, W. T. Kerr, V. S. Haynes, A. L. Yuille, J. Xie, Y. N. Wu, J. A. Brown, and M. S. Cohen. Non-negative matrix factorization of multimodal mri, fmri and phenotypic data reveals differential changes in default mode subnetworks in adhd. *NeuroImage*, 2013.
- [3] L. D. Baker and A. K. McCallum. Distributional clustering of words for text classification. In *Proceedings of the 21st annual*

international ACM SIGIR conference on Research and development in information retrieval, pages 96–103. ACM, 1998.

- [4] D. Cai, X. He, J. Han, and T. S. Huang. Graph regularized nonnegative matrix factorization for data representation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(8):1548–1560, 2011.
- [5] S. Cui, R. Manica, R. F. Tabor, and D. Y. Chan. Interpreting atomic force microscopy measurements of hydrodynamic and surface forces with nonlinear parametric estimation. *Review of Scientific Instruments*, 83(10):103702, 2012.
- [6] S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *JASIS*, 41(6):391–407, 1990.
- [7] G. W. Furnas, S. Deerwester, S. T. Dumais, T. K. Landauer, R. A. Harshman, L. A. Streeter, and K. E. Lochbaum. Information retrieval using a singular value decomposition model of latent semantic structure. In *Proceedings of the 11th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 465–480. ACM, 1988.
- [8] P. X. Gao. Facial age estimation using clustered multi-task support vector regression machine. In *Pattern Recognition (ICPR), 2012 21st International Conference on*, pages 541–544. IEEE, 2012.
- [9] R. He, W.-S. Zheng, and B.-G. Hu. Maximum correntropy criterion for robust face recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 33(8):1561–1576, 2011.
- [10] H. W. Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- [11] L. Li, J. Yang, K. Zhao, Y. Xu, H. Zhang, and Z. Fan. Graph regularized non-negative matrix factorization by maximizing correntropy. *arXiv preprint arXiv:1405.2246*, 2014.
- [12] C.-J. Lin. Projected gradient methods for nonnegative matrix factorization. *Neural computation*, 19(10):2756–2779, 2007.
- [13] X. Liu, Y. Gong, W. Xu, and S. Zhu. Document clustering with cluster refinement and model selection capabilities. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 191–198. ACM, 2002.
- [14] Z. Qin and C. R. Shelton. Improving multi-target tracking via social grouping. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 1972–1978. IEEE, 2012.
- [15] F. Shahnaz, M. W. Berry, V. P. Pauca, and R. J. Plemmons. Document clustering using nonnegative matrix factorization. *Information Processing & Management*, 42(2):373–386, 2006.
- [16] J. Shen and S.-C. S. Cheung. Layer depth denoising and completion for structured-light rgb-d cameras. In *Computer Vision and Pattern Recognition (CVPR), 2013 IEEE Conference on*, pages 1187–1194. IEEE, 2013.
- [17] J. Shen, P.-C. Su, S.-c. S. Cheung, and J. Zhao. Virtual mirror rendering with stationary rgb-d cameras and stored 3-d background. *Image Processing, IEEE Transactions on*, 22(9):3433–3448, 2013.
- [18] H. Song, X. Li, and P. Wang. Image annotation refinement using dynamic weighted voting based on mutual information. *Journal of Software (1796217X)*, 6(11), 2011.
- [19] Q. Sun, F. Hu, and Q. Hao. Mobile target scenario recognition via low-cost pyroelectric sensing system: Toward a context-enhanced accurate identification. 2013.
- [20] Q. Sun, R. Ma, Q. Hao, and F. Hu. Space encoding based human activity modeling and situation perception. In *Cognitive Methods in Situation Awareness and Decision Support (CogSIMA), 2013 IEEE International Multi-Disciplinary Conference on*, pages 183–186. IEEE, 2013.
- [21] Q. Sun, P. Wu, Y. Wu, M. Guo, and J. Lu. Unsupervised multi-level non-negative matrix factorization model: Binary data case. *Journal of Information Security*, 3(4), 2012.
- [22] T. Sun, S. Ding, and Z. Ren. Novel image recognition based on subspace and sift. *Journal of Software (1796217X)*, 8(5), 2013.
- [23] J. J.-Y. Wang, X. Wang, and X. Gao. Non-negative matrix factorization by maximizing correntropy for cancer clustering. *BMC bioinformatics*, 14(1):107, 2013.
- [24] J.-Y. Wang, I. Almasri, and X. Gao. Adaptive graph regularized nonnegative matrix factorization via feature selection. In *Pattern Recognition (ICPR), 2012 21st International Conference on*, pages 963–966, 2012.
- [25] Y. Wang, W. Jiang, and G. Agrawal. Scimate: A novel mapreduce-like framework for multiple scientific data formats. In *Cluster, Cloud and Grid Computing (CCGrid), 2012 12th IEEE/ACM International Symposium on*, pages 443–450. IEEE, 2012.
- [26] L. Xu, Z. Zhan, S. Xu, and K. Ye. Cross-layer detection of malicious websites. In *Proceedings of the third ACM conference on Data and application security and privacy*, pages 141–152. ACM, 2013.
- [27] W. Xu, X. Liu, and Y. Gong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 267–273. ACM, 2003.
- [28] J. Yang and Z. Fei. Broadcasting with prediction and selective forwarding in vehicular networks. *International Journal of Distributed Sensor Networks*, 2013.
- [29] J. Yang, Y. Wang, K. Hua, and W. Wang. Fairness based dynamic channel allocation in wireless mesh networks. In *Computing, Networking and Communications (ICNC), 2014 International Conference on*, pages 556–560. IEEE, 2014.
- [30] Z. Yu, O. C. Au, R. Zou, W. Yu, and J. Tian. An adaptive unsupervised approach toward pixel clustering and color image segmentation. *Pattern Recognition*, 43(5):1889–1906, 2010.
- [31] H. Zhang, Z. Zhang, H. Dai, R. Yin, and X. Chen. Distributed spectrum-aware clustering in cognitive radio sensor networks. In *Global Telecommunications Conference (GLOBECOM 2011), 2011 IEEE*, pages 1–6. IEEE, 2011.
- [32] Y. Zhou, L. Li, and H. Zhang. Adaptive learning of region-based pls model for total scene annotation. *arXiv preprint arXiv:1311.5590*, 2013.

- [33] Y. Zhou, L. Li, T. Zhao, and H. Zhang. Region-based high-level semantics extraction with cedd. In *Network Infrastructure and Digital Content, 2010 2nd IEEE International Conference on*, pages 404–408. IEEE, 2010.